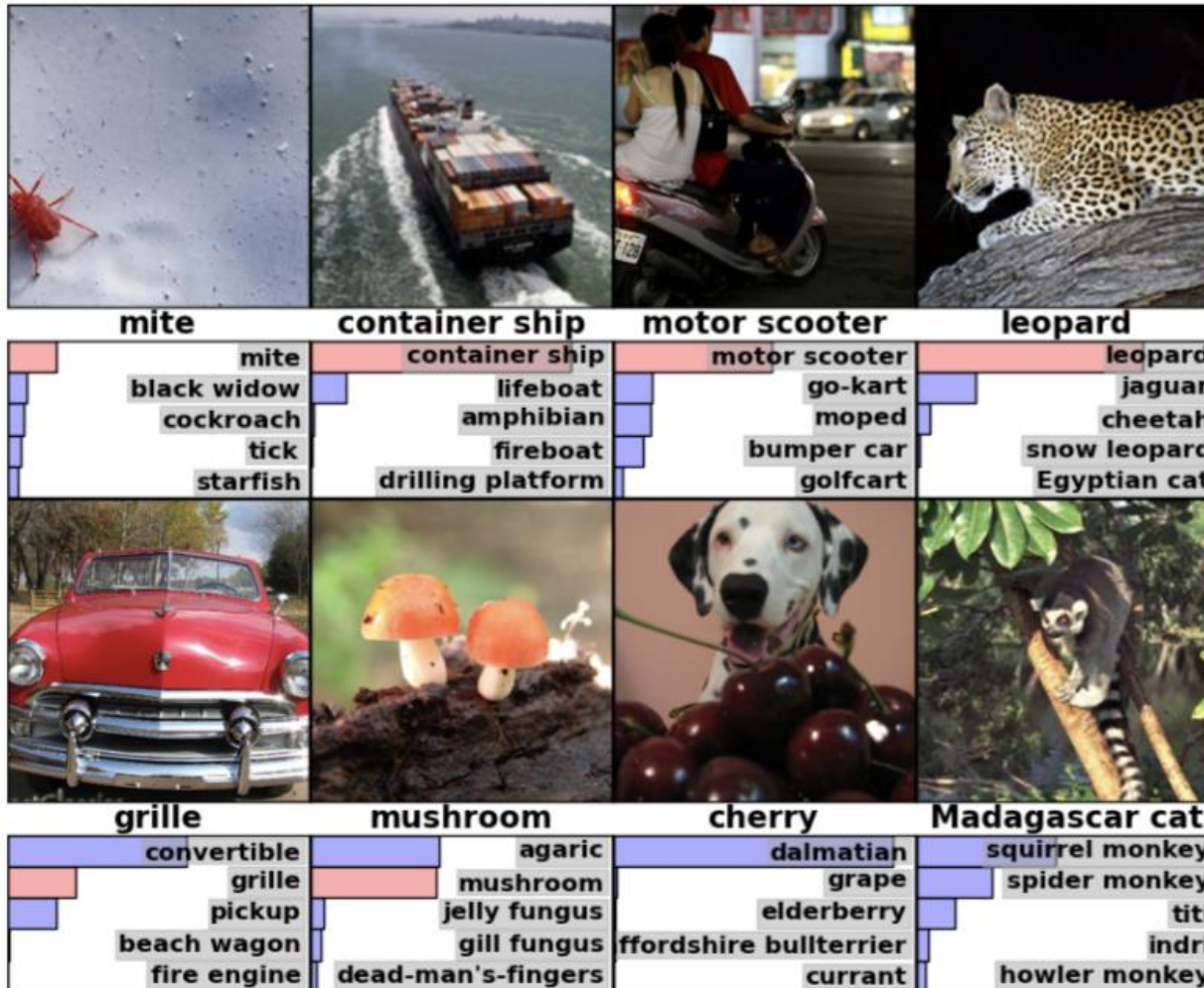


# Adversarial Examples in Machine Learning

---

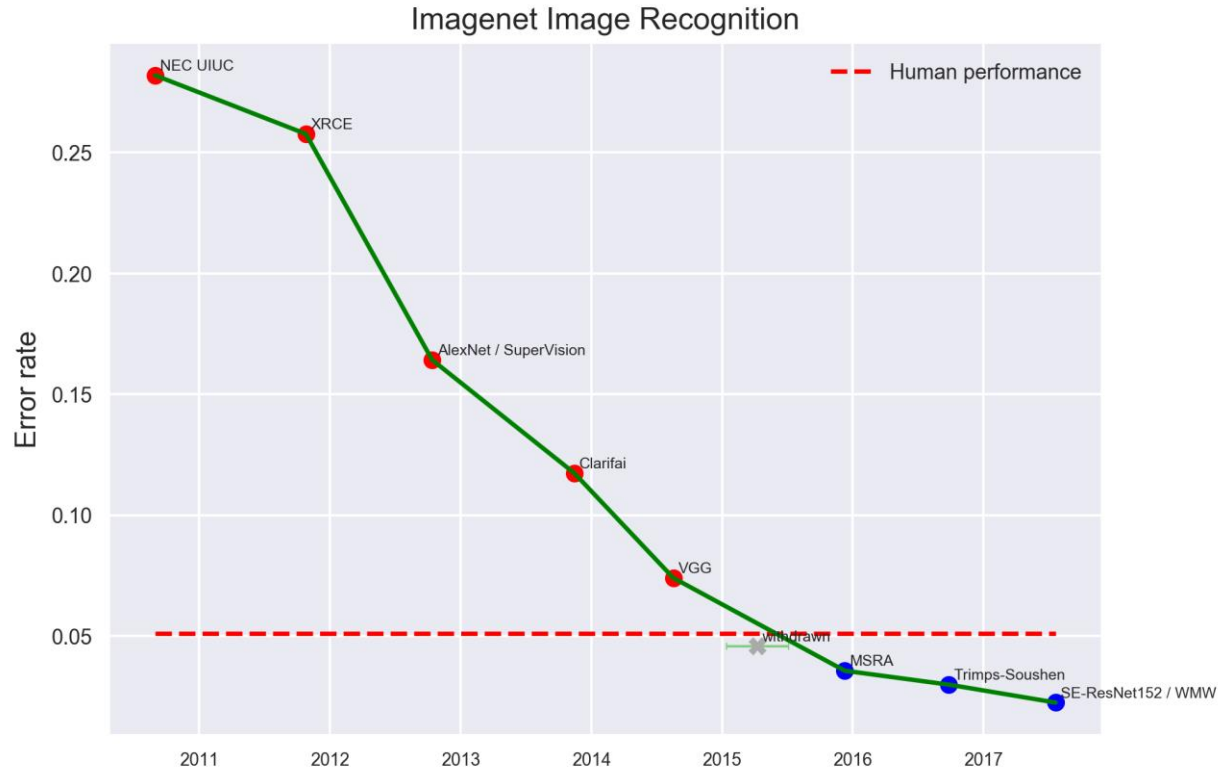
**Robin Jia**  
USC CSCI 467, Spring 2023  
April 18, 2023

# Previously: Image classification

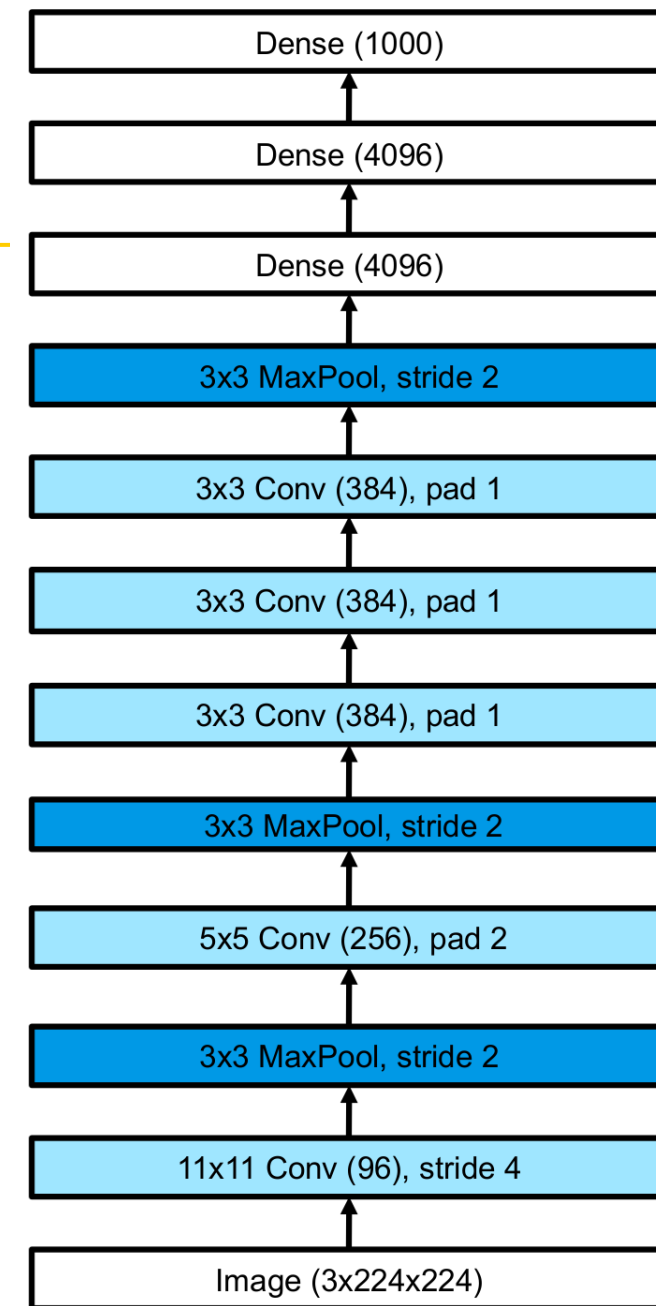


- ImageNet dataset: 14M images, 1000 labels
- **CNNs do very well at these tasks!**

# Previously: ImageNet Progress



- 2012: AlexNet wins ImageNet challenge, marks start of deep learning era (**and is a convolutional neural network**)
- 2016: Machine learning surpasses human accuracy



# Today: A “Reality Check”

---

- Do models really “see” images the way humans do?
- Are models learning shortcuts rather than actually solving the task?

**Adversarial Examples**

**Spurious Correlations**



# Adversarial Examples

- **Adversarial examples:** Examples crafted by an **adversary** (attacker) to cause a desired behavior by a machine learning model
  - Can exist despite high average accuracy



Panda  
58% confidence

+ .007 ×



Nematode  
8% confidence

=



Gibbon  
99% confidence



■ classified as turtle    ■ classified as rifle  
■ classified as other

Goodfellow et al. "Explaining and Harnessing Adversarial Examples." ICLR 2015.

Sharif et al. "A General Framework for Adversarial Examples with Objectives." ACM TOPS 2019.

Athalye et al. "Synthesizing Robust Adversarial Examples." ICML 2019.

# Why do we care?

---



## Security

- Fooling facial recognition systems
- Vulnerabilities of safety-critical systems (e.g. self-driving cars)
- Bypassing content moderation or spam detection
- Hacking ranking algorithms (search engine optimization)



## Interpretability

- Do models work the way we think they do?
- Understand model weaknesses so we can patch them
- Understand when models might not be reliable

# The rules of the game

---

## Defining the **threat model**

1. **Attack vector:** What can the adversary do?
2. **Adversary's knowledge:** What does the adversary know?
3. **Adversary's goal:** What does the adversary want to achieve?



# Attack vectors



- Apply a perturbation to input (Constrained attack)



Panda  
58% confidence

+ .007 ×



Nematode  
8% confidence

=



Gibbon  
99% confidence



# Attack vectors



- Apply a perturbation to input (Constrained attack)
- Completely change the input (Unconstrained attack)
- Add bad training data (Data poisoning)



■ classified as turtle    ■ classified as rifle  
■ classified as other

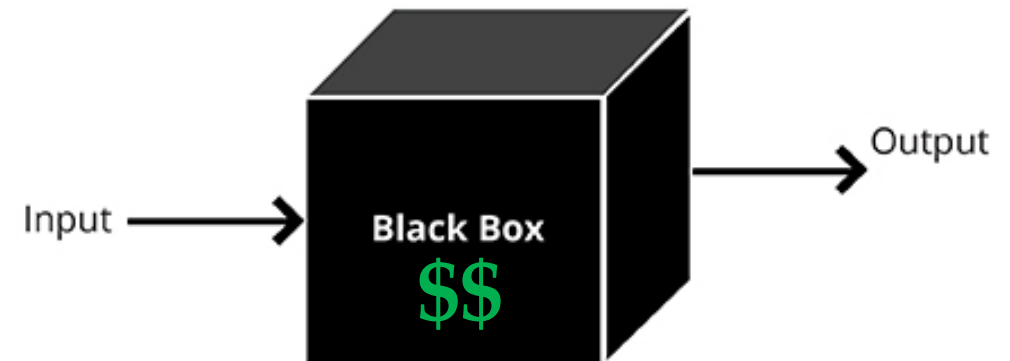
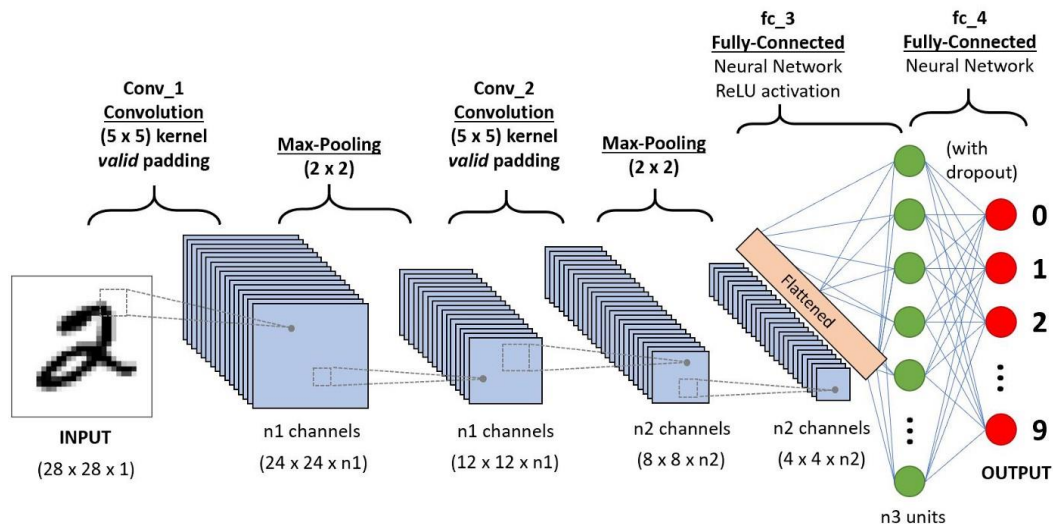
# Adversary's knowledge



**White-box:** Has access to model and all internals (e.g., has model parameters and code)

**Black-box:** Has access to model only via queries

- May also have a query budget

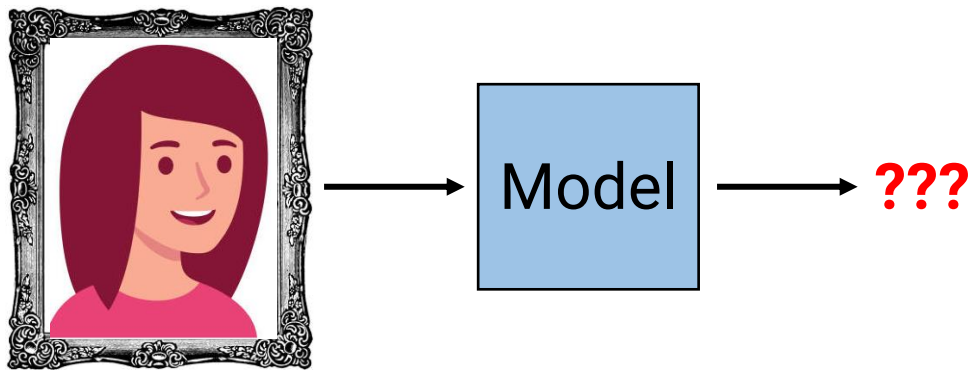


# Adversary's goal



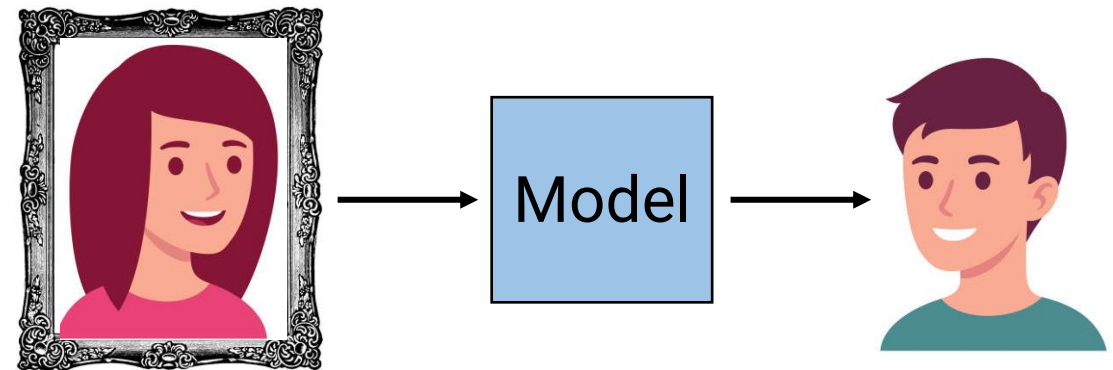
## Undirected: Cause any error

- Facial recognition: Avoid being detected as yourself



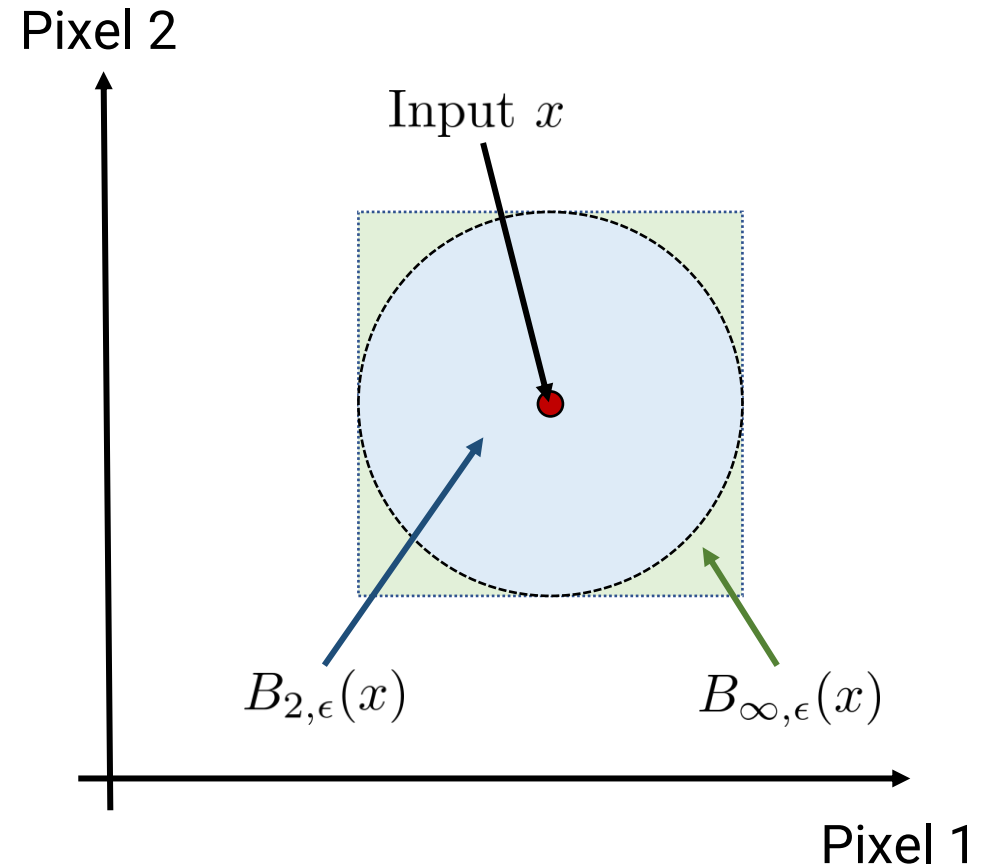
## Directed: Cause a specific (wrong) prediction

- Facial recognition: Appear to be some other specific person



# Adversarial perturbations for images

- Informal attack vector: Make imperceptible change to image
- How to formalize?
  - Make new image  $x'$  very close to  $x$  in pixel space
    - L2 norm:  $\|x_i - x\|_2 = \sqrt{\sum_{i=1}^d (x'_i - x_i)^2}$
    - L-infinity norm:  $\|x_i - x\|_\infty = \max_i |x'_i - x_i|$
  - Constrain norm of difference to be small, e.g.  $\|x' - x\|_\infty \leq \epsilon$ 
    - Equivalently,  $x' \in B_{\infty, \epsilon}(x)$
    - Each pixel can change by  $\epsilon$



# Adversarial perturbations for images

- The rules of the game
  - Attack vector: Given test example  $x$ , replace with any  $x' \in B_{\infty, \epsilon}(x)$
  - Informally: Attacker can change brightness of each pixel by at most  $\epsilon$
  - Knowledge: White box
  - Goal: Undirected (could also be directed for multi-class)



Panda  
58% confidence

+ .007 ×



Nematode  
8% confidence

=



Gibbon  
99% confidence

# Attacking a classifier

---

- Problem statement for attacker
  - Binary classification, model predicts  $\text{sign}(f(x; \theta))$
  - Given: Image  $x$ , label  $y$ , model parameters  $\theta$
  - Return:  $x' \in B_{\infty, \epsilon}(x)$  such that  $\text{loss}(x', y; \theta)$  is maximized

# Attacking a classifier

- Approximate solution (“Fast Gradient Sign Method” or FGSM )

- Let  $z = x' - x$

- Idea: Approximate  $f$  locally with a linear model

$$f(x'; \theta) \approx f(x; \theta) + \underbrace{\nabla_x f(x)^\top (x' - x)}_{\text{Original prediction}} = f(x; \theta) + \underbrace{\nabla_x f(x)^\top z}_{\text{Adversary controls}}$$

Gradient with respect to  $\mathbf{x}$  (not the parameters!)

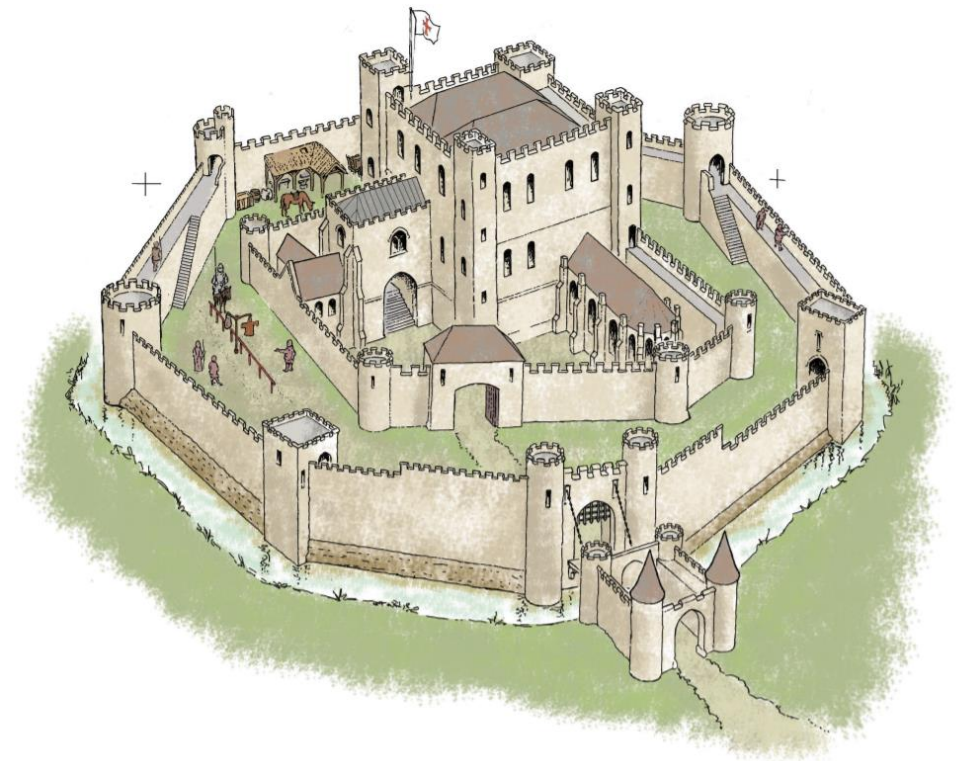
- To increase  $f$ , add  $\varepsilon$  when gradient  $> 0$ , subtract  $\varepsilon$  when gradient  $< 0$
  - Do the reverse if adversary wants to decrease  $f$

|                               |                |                |   |                |   |
|-------------------------------|----------------|----------------|---|----------------|---|
| $\nabla_x f(x)$               | 1.2            | -2.8           | 0 | 2.3            |   |
| $z$ to <b>increase</b> $f(x)$ | $\varepsilon$  | $-\varepsilon$ | 0 | $\varepsilon$  | (Adversary makes model predict $y=+1$ ) |
| $z$ to <b>decrease</b> $f(x)$ | $-\varepsilon$ | $\varepsilon$  | 0 | $-\varepsilon$ | (Adversary makes model predict $y=-1$ ) |

# Defending against adversarial perturbations

---

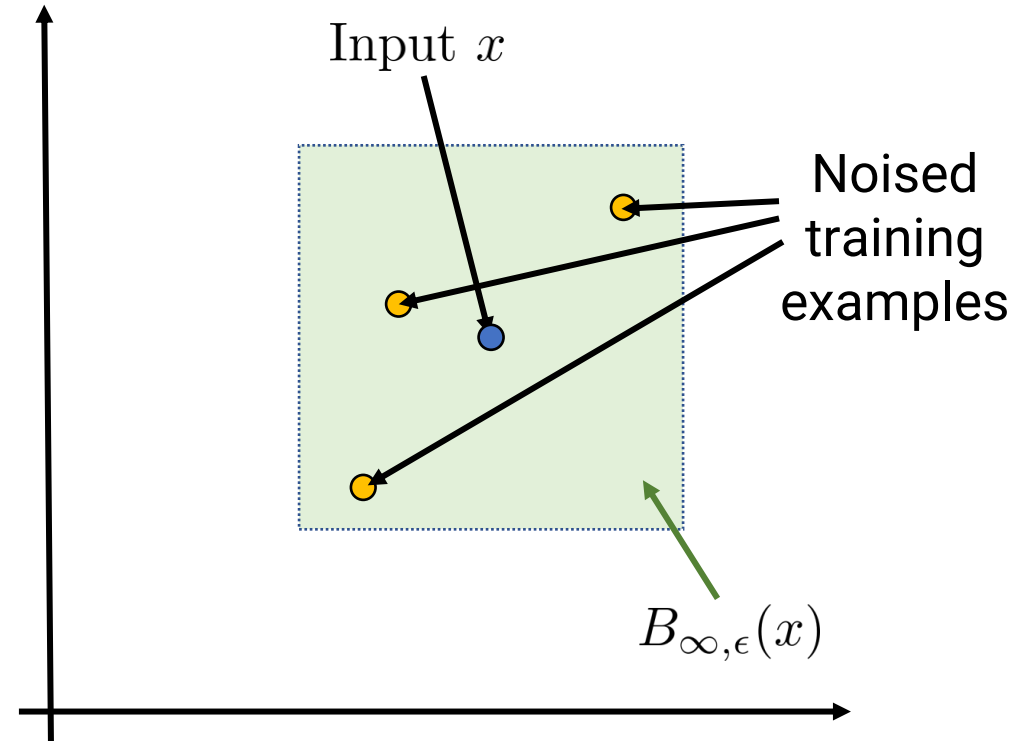
- Problem statement for defender
  - Given: Dataset  $D$  and **known threat model**
    - i.e. Assume you know the norm and perturbation radius  $\epsilon$
  - Return: Model parameters  $\theta$  such that attacker cannot succeed
- Adversary has second player advantage!
  - First, you train the model
  - Then the adversary gets to attack it



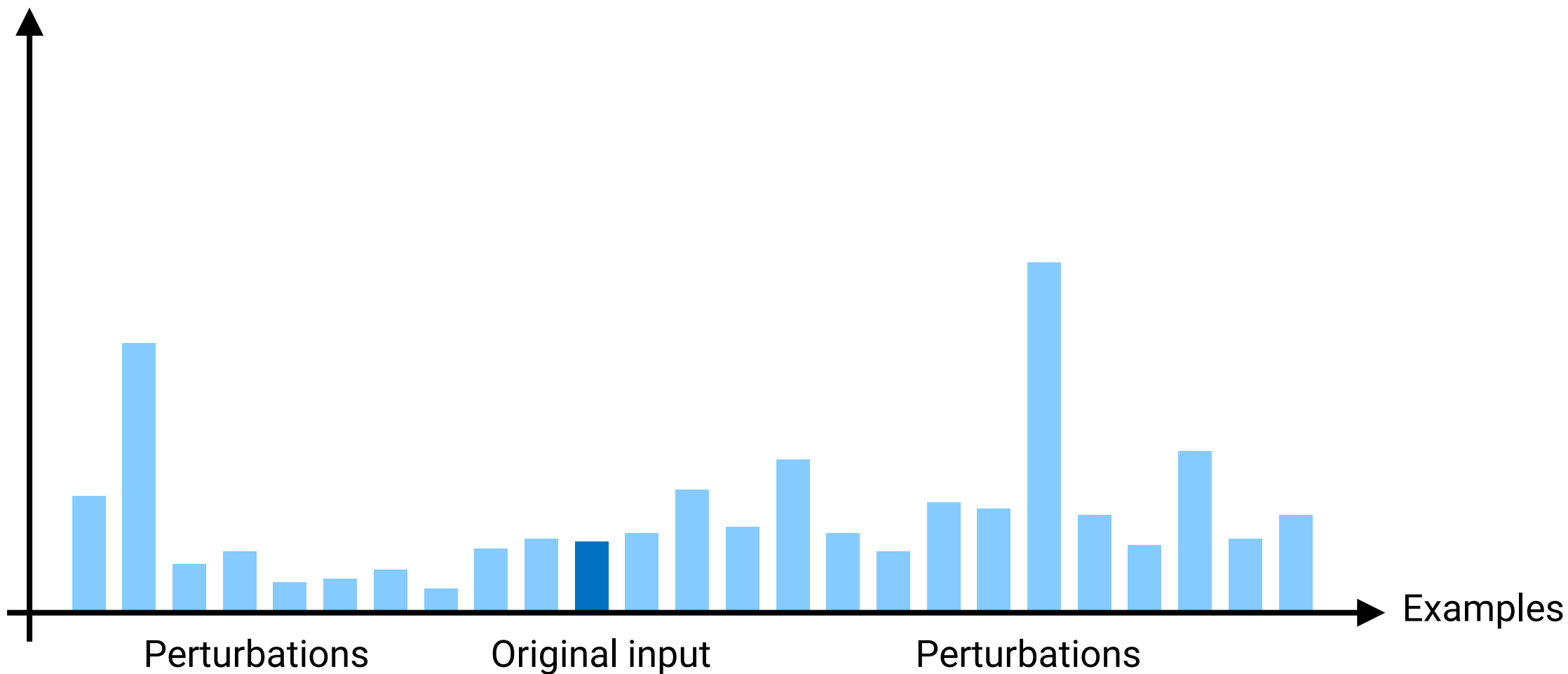


# A naïve defense strategy

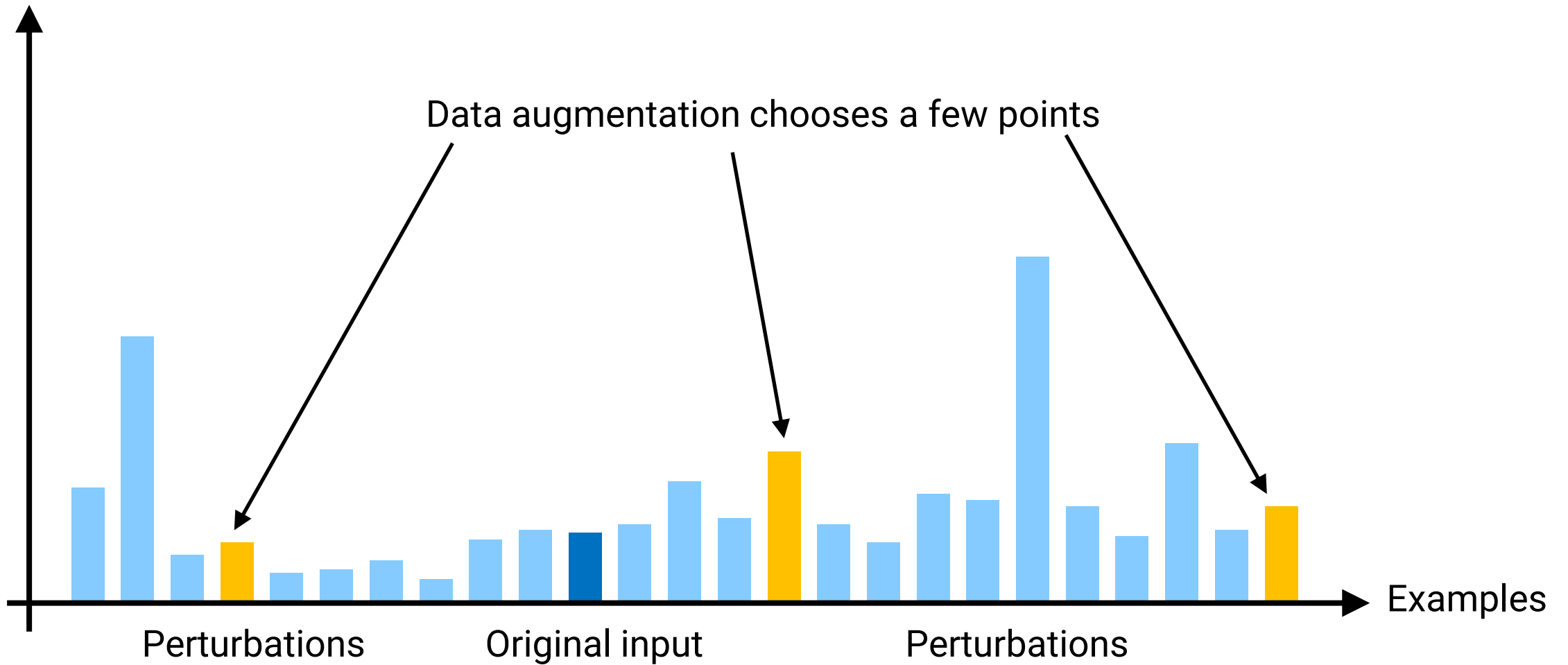
- **Data augmentation:** Automatically generate additional training examples based on your current data
  - Often a good strategy in general, but not here...
- Random data augmentation
  - Randomly add noise to training examples  $x$  within  $B_{\infty, \epsilon}(x)$
  - Train on this augmented data
- Problem: Adversary is choosing worst-case perturbation, may be much worse than random perturbation!



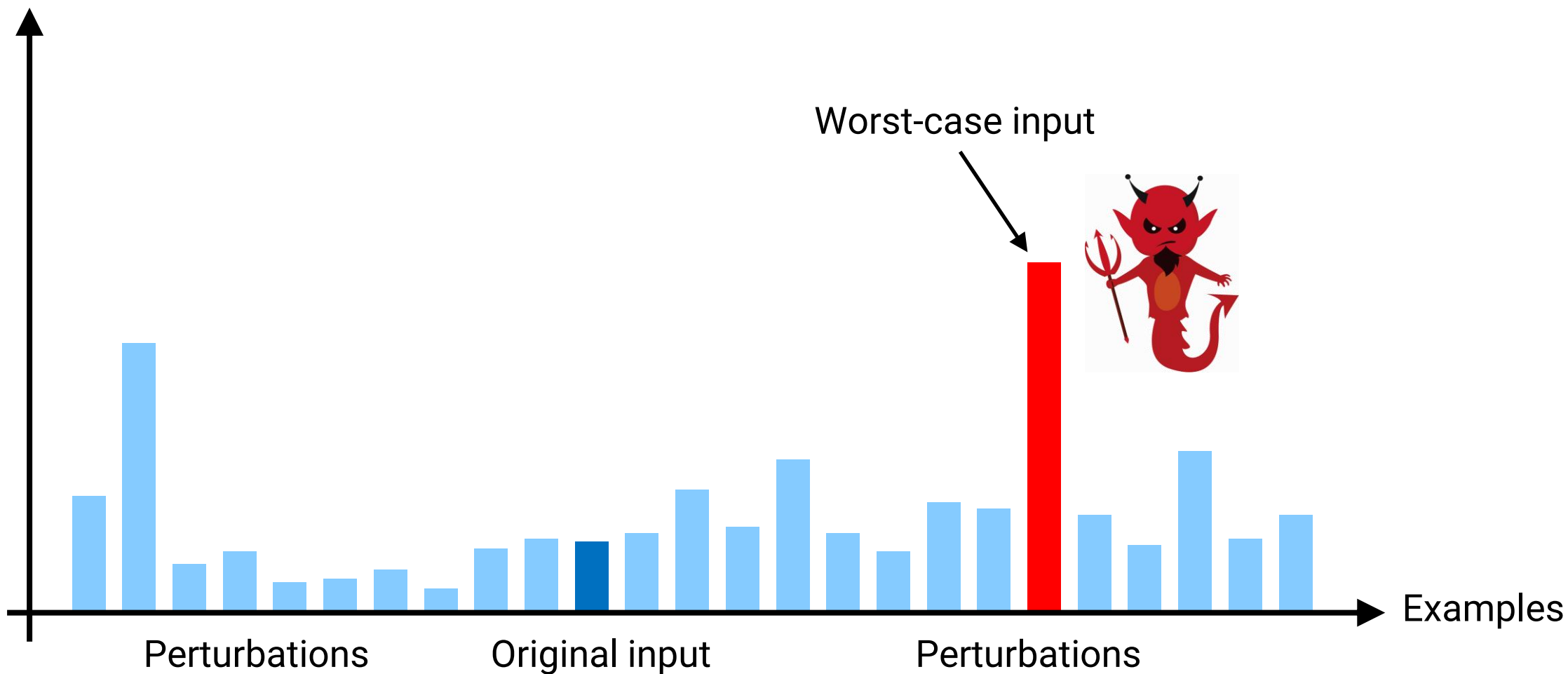
Loss (lower = better)



Loss (lower = better)

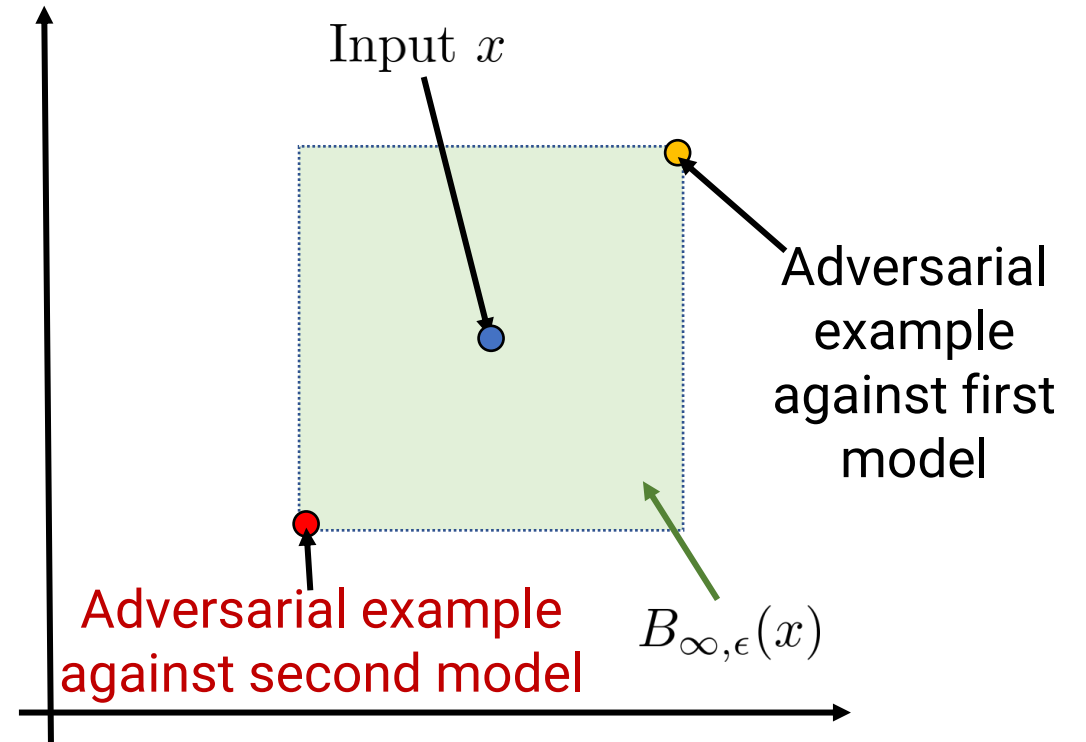


Loss (lower = better)



# Another naïve defense strategy

- “Adversarial data augmentation”
  - Train model normally
  - Generate adversarial examples for this model
  - Add these to training data and retrain
- Flaw: At test time, adversary can perturb in a different way!



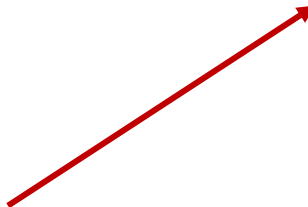
# Anticipating the adversary

---


- Normal training loss function: 
$$\min_{\theta} \sum_{(x,y) \in D} \text{loss}(x, y; \theta)$$

- What we want to optimize instead: 
$$\min_{\theta} \sum_{(x,y) \in D} \max_{x' \in B_{\epsilon}(x)} \text{loss}(x', y; \theta)$$

Choose the parameter that minimizes training loss...



On the perturbation that the optimal adversary would choose **against this model!**



# Adversarial training

---

- How can we optimize  $\min_{\theta} \sum_{(x,y) \in D} \max_{x' \in B_{\epsilon}(x)} \ell(y \cdot f(x'; \theta))$  ?
- Run an attack algorithm A (e.g., FGSM) **against current model** to generate  $x' = A(x, y; \theta)$
- Plug it in:  $\min_{\theta} \sum_{(x,y) \in D} \ell(y \cdot f(\underbrace{A(x, y; \theta)}_{\text{Adversarial example for current model}}); \theta)$
- Implementation: Every time you want to do a gradient step, first run the attack, then do gradient step on the adversarial example

# NLP: Adversarial Unicode attacks

---

- Images: We could have some **actually imperceptible** perturbations
- Text equivalent: Unicode characters that look like ASCII characters

## I. INTRODUCTION

Do `x` and `х` look the same to you? They may look identical to humans, but not to most natural-language processing systems. How many characters are in the string `"123"`? If you guessed 100, you're correct. The first example contains the `Latin character x` and the `Cyrillic character h`, which are typically rendered the same way. The second example contains `97 zero-width non-joiners1` following the visible characters.

<sup>1</sup>Unicode character U+200C



# NLP: Typo-based attacks

- Adversarially chosen typos can also cause misclassification
- Think about an RNN or Transformer
  - Input is a set of word vectors
  - Add a typo = completely different word vector for that word!

| Alteration | Movie Review  | Label |
|------------|---|-------|
| Original   | A triumph, relentless and beautiful in its downbeat darkness        | +     |
| Swap       | A triumph, relentless and <b>beuatiful</b> in its downbeat darkness | -     |
| Drop       | A triumph, relentless and beautiful in its <b>dwnbeat</b> darkness  | -     |
| + Defense  | A triumph, relentless and <b>beautiful</b> in its downbeat darkness | +     |
| + Defense  | A triumph, relentless and beautiful in its <b>downbeat</b> darkness | +     |

# NLP: Meaning preserving attacks

- Can keep meaning the same (e.g. *“What has”* -> *“What’s”*)
- Security case
  - Alter model prediction while maintaining equivalent meaning to a reader
    - SEO, Plagiarism detection
- Interpretability case
  - Surprising if model succeeds on one input but fails on another that people would think of as equivalent

In the United States especially, several high-profile cases such as Debra LaFave, Pamela Rogers, and Mary Kay Letourneau have caused increased scrutiny on teacher misconduct.

(a) Input Paragraph

**Q:** What has been the result of this publicity?  
**A:** increased scrutiny on teacher misconduct

(b) Original Question and Answer

**Q:** What **haL** been the result of this publicity?  
**A:** **teacher misconduct**

(c) Adversarial Q & A (Ebrahimi et al., 2018)

**Q:** **What’s** been the result of this publicity?  
**A:** **teacher misconduct**

(d) **Semantically Equivalent Adversary**

# Summary: Adversarial Examples

---

- White-box attack strategy (**Fast Gradient Sign Method**)
  - Optimal for linear model (Homework!)
  - Approximate for neural model
- Training-time defense (**Adversarial Training w/ FGSM**)
  - Guards against optimal attack for linear model (Homework!)
  - **Guards against approximate attack for neural model**
- Most famous in images, but can occur in any modality
- Lots of research on more sophisticated attacks/defenses, what this means for deployed models, etc.

# Announcements

---

- Homework 4 released, due Thursday, April 27
  - You should be ready to do everything after today's lecture
- Final Exam
  - Thursday, May 4, 2:00-4:00pm (2 hours)
  - Allowed **two** (double-sided) pages of notes
  - Cumulative exam, somewhat more weight on material after midterm
  - Similar in style to midterm exam
- Last section this Friday
  - Review of course material

# Today: A “Reality Check”

---

- Do models really “see” images the way humans do?
- Are models learning shortcuts rather than actually solving the task?

**Adversarial Examples**

**Spurious Correlations**



# Previously: Machine learning is a tornado

---

- ...it picks up everything in its path
- Data has all sorts of associations we may not want to model



# Some pictures of wolves

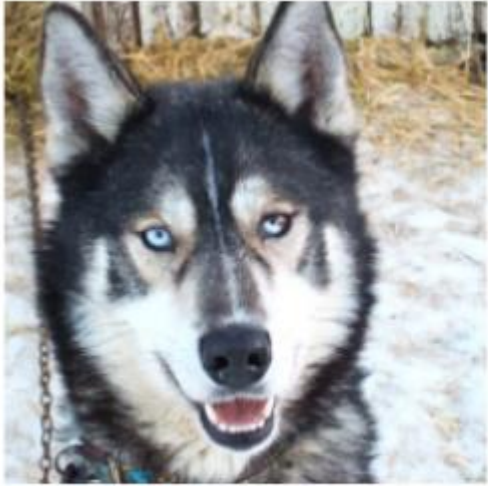
---



What do these have in common...?

# What does the model learn?

---



(a) Husky classified as wolf



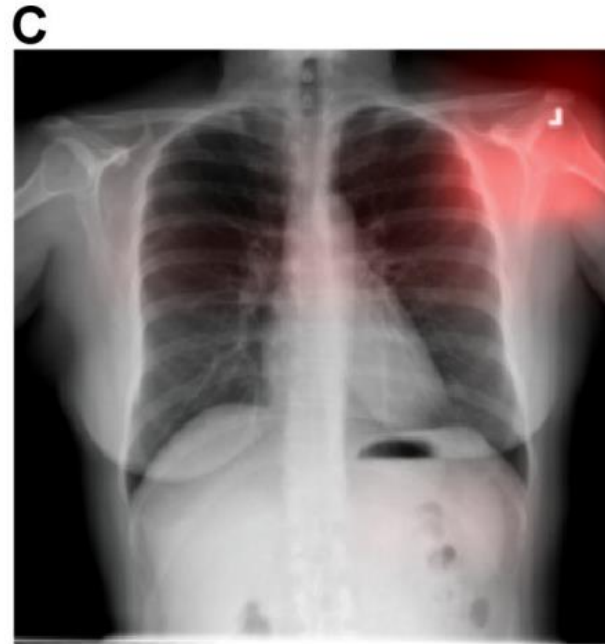
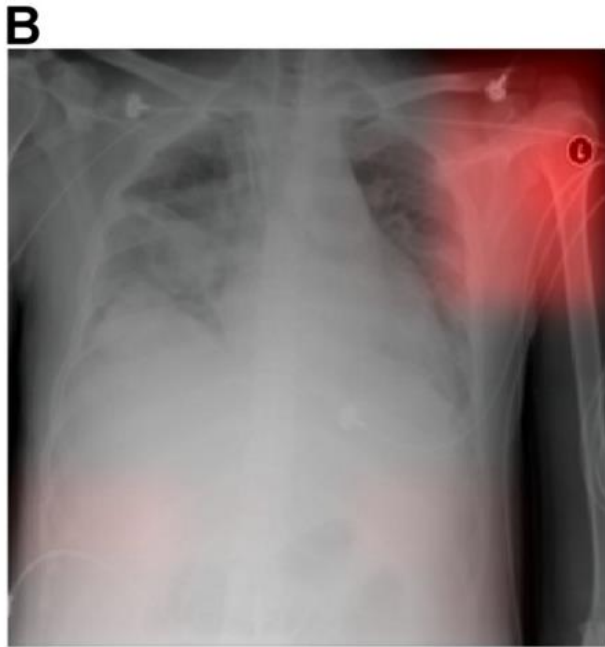
(b) Explanation

- Model misclassifies husky as a wolf
- Why? Model sees snow and associates it with huskies
- This is a **spurious correlation**
  - Model is just trying to associate input features with label
  - Snow is correlated with “wolf” label, so model learns this
  - But this is **spurious**—not part of the actual task



# Spurious correlations in medicine

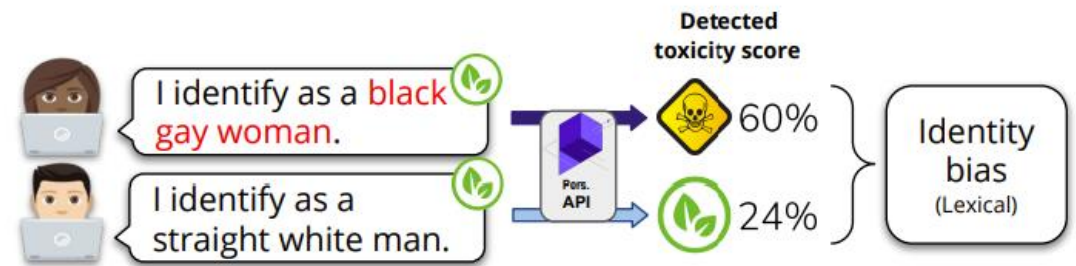
---



- Task: Detecting pneumonia from chest X-ray
- Spurious correlation: Metallic token radiology technicians place on patient
  - Different hospitals do this differently
  - Different hospitals have different pneumonia prevalence
- Result: Model relies heavily on these hospital-specific tokens!

# Spurious correlations in NLP

- Hate speech detection: Identity mentions lead to model predicting text as toxic
  - Spurious correlation: Hateful speech directed at specific groups often names those groups
- Sentiment analysis: Some names associated with positive/negative sentiment



| Sentence                  | Toxicity | Sentiment |
|---------------------------|----------|-----------|
| I hate Justin Timberlake. | 0.90     | -0.30     |
| I hate Katy Perry.        | 0.80     | -0.10     |
| I hate Taylor Swift.      | 0.74     | -0.40     |
| I hate Rihanna.           | 0.69     | -0.60     |

# Spurious correlations and generalization

---

## Common training examples

Waterbirds

y: waterbird  
a: water  
background



y: landbird  
a: land  
background



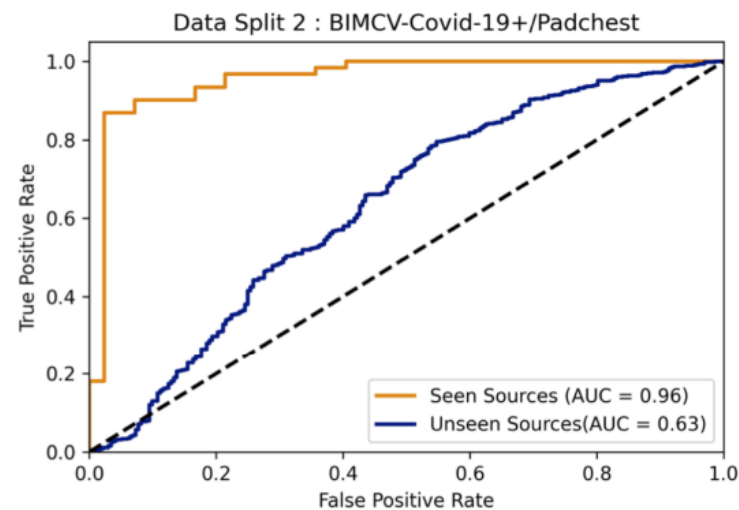
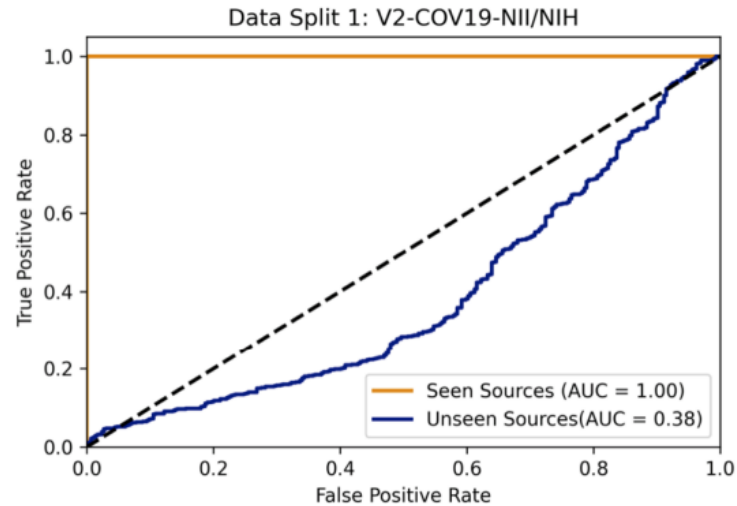
## Test examples

y: waterbird  
a: land  
background



- Task: Identifying bird species
- Spurious correlation: Waterbirds tend to be pictured over water
- Generalization challenge: Cannot identify ducks on land!
  - In general: Overreliance on spurious correlations means your model will perform poorly in scenarios where the correlation no longer holds

# Spurious correlations and generalization



- Task: Detecting pneumonia from chest X-ray (again) in COVID patients
- Compared two settings
  - Seen sources: Train and test on same data sources
  - Unseen sources: Train and test on different data sources (datasets from 3 different countries)
- Model can be very good on seen sources but worse than random on unseen sources!
  - Likely learns source-specific correlations
  - Similar to HW1 and author identification

# Avoiding overreliance on spurious correlation

---

- Lots of research, but no guaranteed solutions
- Diversifying dataset often helps
- General recommendation: Evaluate **out-of-distribution generalization**
  - Go beyond the hospitals you trained on
  - Find pictures of wolves in atypical backgrounds
- Practice caution: Don't assume model will generalize without measuring first



# Conclusion

- Supervised learning learns patterns from a training dataset
- Things can go very wrong when the test data deviates in some way from the training data
  - Addition of adversarial perturbations
  - New data that breaks spurious correlations in the training data
- Careful evaluation is critical to identify these issues

